# New C++ paradigms applied to the implementation of efficient arbitrary-rank tensors for high-performance computing

Alejandro M. Aragón

**Abstract**

This article discusses an efficient implementation of tensors of arbitrary rank that can be used as a basic building block for high-performance computing through the use of some of the C++ paradigms introduced by the new C++11 requirements of the standard library. By using templates, an extra high-level layer can be added to the C++ language when dealing with algebraic objects and their operations, without compromising performance.

## 1 Introduction

Developed in the early 1950s, Fortran has been regarded as the programming language of choice for scientific computing. Nevertheless, scientists are increasingly adopting other languages for high-performance computing motivated by the use of higher level paradigms, object orientation, data abstraction and generic programming. One of these languages is C++, developed in the late 1970s to enhance the C language [1] and standardized in 1998 [2]. Contrary to Fortran, the C++ programming language does not provide built-in support for matrices. As a result, a programmer is forced to either write custom classes that represent the abstraction of algebraic objects, or to use an external library (e.g. PETSc [3], Blitz++ [4], or MTL [5, 6]). In the former case, a naive implementation could drastically deteriorate the performance and lead the user to discard completely the programming language. In the latter, the developer increases the dependency of the resulting code, for the user has to build the library used and its dependents. The addition of a library to handle algebraic objects also raises the important question: if an external library is needed to cope with the absence of such an elemental component in scientific computing, how many other libraries need to be added to the dependency list for a final software design? In this paper the first approach is adopted to craft an efficient implementation of an arbitrary-rank tensor by using the paradigms introduced by the new C++ standard. The methodology described herein is also available as an open-source header-only library for further reference [7].

1

Even though some of the components discussed here are known to any programmer who has been exposed to the C++ programming language, other idioms require advanced knowledge, specially when dealing with metaprogramming, i.e., the generation of code at compilation time. Programmers usually think of a program as a set of instructions that do work at run time. However, some of the idioms used here will deal with compile time programming. For the more advanced topics, the reader is referred to [8, 9] and the references therein. The recently published C++ ISO Standard [10] aims at enhancing the programming language thorough the set of `C++11` requirements. The introduction of *variadic templates* [11, 12] brings unprecedented potential for template metaprogramming. Traditionally, C++ templates contained a fixed number of template parameters, but with the addition of variadic templates this restriction is eliminated. Templates are therefore the main language component used in this paper.

The Basic Linear Algebra Subprograms (BLAS) library [13, 14, 15] is usually the software of choice to carry out algebraic operations with vectors and matrices. A highly optimized implementation of the BLAS library can be found in any high-performance scientific computing environment. However, the library interface is not straightforward to use, as the fourteen-parameter function to carry out multiplication between matrices in Listing 1 shows.

Listing 1: Level 3 BLAS cblas_dgemm function

```
void cblas_dgemm(const enum CBLAS_ORDER Order,
                 const enum CBLAS_TRANSPOSE TransA,
                 const enum CBLAS_TRANSPOSE TransB,
                 const int M, const int N, const int K,
                 const double alpha, const double *A,  const int
                     lda,
                 const double *B, const int ldb, const double beta,
                 double *C, const int ldc);
```

Therefore, the BLAS interface is prone to user error, and thus a more friendly interface is desirable even if the BLAS library is still used at the backend. One such project is the uBLAS library [16], that is part of the BOOST set of C++ libraries. When dealing with custom classes that represent algebraic objects, an amiable syntax can be created by using *operator overloading*. That is, the user can write `C=A*B` to obtain the product between two matrices A and B, even if the implementation calls the aforementioned BLAS function shown in Listing 1. Through *expression templates* [8], the compiler can determine the precise function call at compilation time without adding extra running time overhead. But probably the most powerful feature described in this manuscript is the capability to tailor the syntax for specific operations. By using expression templates, it is possible to match complicated expressions in order to use the complex higher-level routines from the BLAS library with mathematical-like syntax. It is worth mentioning that both the Blitz++ [4] and the MTL [6] C++ libraries were built taking advantage of the template facility. While the former was built entirely on the concept of expression templates, the latter uses generic

programming techniques for achieving high performance.

This article is organized as follows: Some of the C++ features introduced by the new C++11 requirements of the standard library are summarized in Section 2. Section 3 shows the use of new C++ paradigms in the implementation of an **Array** class template that can be used for the definition of a tensor of any rank, i.e., a template class that can define vectors, matrices, and arbitrary-rank tensors. By using operator overloading, Section 4 sets forth the syntax that provides the user with an accessible mathematical notation without undermining performance. It presents the implementation of the classes and operators that define the language that allows us to use the algebraic objects in mathematical expressions. Performance tests and usage are presented in Section 5.

## 2 New C++11 language features

Most of the new features from the C++11 set of requirements used in this work are illustrated in the following code snippet below, whose details can be found in the C++ ISO Standard [10].

```
template <typename... Types> void foo (Types... args);

void bar() {}
template <typename T, typename... Types> void bar(T u, T v, Types
    ... args) {
  static_assert(sizeof...(args) % 2 == 0, "Odd number of
      parameters");
  cout<<u<<"-"<<v<<endl; bar(args...);
}

template <typename... Types> class ClassA;
template <typename T, typename... Types> class ClassB;

template <typename... Types>
using ClassBInt = B <int, Types... >;
```

Prior to C++11, the C++ programming language only supported a fixed number of template parameters for both classes and functions. In the snipped above, Types is a *parameter pack identifier*, and the new syntax for variadic templates uses the ellipsis operator to indicate that functions foo, bar, and classes ClassA, ClassB, can take an arbitrary number of parameters. Notice the different placement of the ellipsis operator with respect to Types in the declarations above. Placing the operator to the left declares the identifier Types as a *template parameter pack*, whereas placing it to the right denotes a *parameter pack expansion*. In the declaration of the functions above, the pack expansion denotes a *function parameter pack*. A pack expansion can be used in the body of a function, e.g., to produce a list of instantiated elements that match their corresponding type in the pack. The function **bar** above uses variadic templates to print ordered pairs given an arbitrary even number of parameters. The definition of the function makes use of recursive calls with a pack

3

expansion, and the recursion terminates with a call to the base case empty function. The **sizeof...** operator, also introduced in the new standard, can be used to produce a compile time constant representing the number of elements in a parameter pack.

Compile time assertions constitute another component, so that the compiler can determine errors at compilation time based on a boolean constant. This feature is also introduced in the new standard, even though the technique is not new and it can be implemented with the previous standard [9]. The compile time assertion is accomplished through the **static_assert** declaration shown in the snippet, which produces a compile error if the function is called with an odd number of parameters.

The last declaration in the code snippet is an *alias template*, and the resulting identifier ClassBInt represents a family of types that take an arbitrary number of template parameters but that have the first parameter as an **int**. All of these additions to the standard library are used in the next section to craft an efficient arbitrary-rank array class template.

## 3   Arbitrary-rank tensor class template

Vectors, matrices and even higher-rank tensors can be obtained from a single **Array** class template with the idioms presented in the previous section. In this way, an array for a specific rank can be determined with a simple alias template declaration:

```cpp
template <int k, typename T> class Array; // array class template
    declaration

template <class T>
using vector_type = array::Array<1,T>; // family of vectors

template <class T>
using matrix_type = array::Array<2,T>; // family of matrices

template <class T>
using tensor_type = array::Array<4,T>; // family of 4th-order
    tensors
```

Yet, the interface of the resulting tensor would depend on its actual rank. For example, a matrix would have functions to return the number of rows and columns, which do not make sense for a vector. Thus, to accommodate the different interfaces depending on the rank, the **Array** class template is a subclass of an **Array_base** class template, displayed in listing 2. The latter stores the state by means of two member variables that are used for the array dimensions and its data. The second class template in Listing 2 is a partial template specialization for two-dimensional arrays. In this case, two member functions are defined to return the number of rows and columns, providing the functionality described above. Thus, through partial template specialization, the developer can define the interface (and state) for tensors of a specific rank.

4

**Listing 2: Array_base class template**

```cpp
template <int k, typename T>
struct Array_base {

  typedef T value_type;
  Array_base() : n_(), data_(NULL) {}
  // ...

protected:
  size_t n_[k];
  value_type* data_;
};

// partial template specialization for a matrix
template <typename T>
struct Array_base<2,T> {

  typedef T value_type;
  Array_base() : n_(), data_(NULL) {}

  size_t rows() const;
  size_t columns() const;
  // ...

protected:
  size_t n_[2];
  value_type* data_;
};
// ... other partial templare specializations
```

The creation of vectors, matrices and other tensors from a single class template is possible due to the introduction of variadic templates [12]. The user can create algebraic objects by calling the same constructor, which takes a variable number of parameters depending on the rank of the object, as follows:

```cpp
vector_type<double> x(3);
matrix_type<float> A(3,3), B(3);
tensor_type<int> II(3,3,3,3), TT(3);
```

The constructor, shown in Listing 3, contains a compile time assertion to ensure that the number of parameters passed is at most equal to the rank of the array. A parameter pack expansion follows and array dimensions are set. Finally, it allocates and initializes all the array elements in a one-dimensional array. The array elements are laid out in memory this way purposely to take advantage of the BLAS routines. The copy constructor, assignment operator and destructor of the **Array** class template follow the traditional C++ idioms.

**Listing 3: Variadic template constructor**

```cpp
// inside Array class template
```

```cpp
typedef Array_base <k,T> base_type;

template <typename... Args>
explicit Array(const Args&... args) : base_type() {

  // check number of parameters
  static_assert(sizeof...(Args) <= k , "Number of arguments
      exceeded rank");
  // unpack parameters
  size_t sizes[] = {args...};
  // assign array sizes
  size_t size = 1;
  for (size_t i=0; i<k; ++i) {
    n_[i] = i < sizeof...(Args) ? sizes[i] : sizes[sizeof...(Args
      ) - 1];
    assert(n_[i] != 0);
    size *= n_[i];
  }
  // allocate and initialize array elements
  data_ = new value_type[size];
  for (size_t i=0; i<size; ++i)
    data_[i] = value_type();
}
```

Element access to the **Array** can be accomplished by overloading **operator()** and **operator[]**. As before, variadic templates are used to provide element access through an overloaded **operator()** that uses variadic template syntax, as illustrated in Listing 4. The function uses again a compile time assertion on the number of parameters passed and calls a private member function to obtain the index in the one-dimensional array that contains the required element. Note however that accessing an element this way requires the computation of an index in the one-dimensional array mentioned above. This function call can be avoided by providing iterators to the class if certain access patterns are known (e.g., iterate over a row or a column of a matrix).

Listing 4: Variadic template **operator()**

```cpp
// inside Array class template
template <typename... Args>
reference_type operator()(Args... params) {
  // check number of parameters
  static_assert(sizeof...(Args) == k , "Arguments-rank mismatch"
      );
  // unpack parameters
  size_t indices[] = {params...};
  // return reference
  return data_[index(indices)];
}
private:
  size_t index(size_t indices[]) const {
```

```
      size_t  i = 0,  s = 1;
      for (int  j=0; j<k; ++j) {
        assert(indices[j] < n_[j]);
        i += s * indices[j];
        s *= n_[j];
      }
      return i;
    }
```

The traditional C-style means of accessing an array element through **operator[]** can also be implemented for the **Array** class template. Because **operator[]** is a unary operator, i.e., an operator that takes a single argument, the technique relays on returning a proxy object from **operator[]** when dealing with arrays for which $k > 1$ [17]. For example, if the array represents a matrix (i.e., $k = 2$), the returned proxy object also implements **operator[]** to finally return an element of the matrix. However, in our **Array** class template some metaprogramming is needed for the array has arbitrary rank. Listing 5 shows the implementation of **operator[]** and the generalized **Array_proxy** class template.

Listing 5: Unary **operator[]**

```
  // inside Array class
  typedef typename Array_proxy<k, Array>::reference_type
      proxy_reference_type;
  typedef typename Array_proxy<k, Array>::value_type
      proxy_value_type;

  proxy_reference_type operator[](size_t i)
  { return proxy_reference_type(*this, i); }

  proxy_value_type operator[](size_t i) const
  { return proxy_value_type(*this, i); }
  // ...
};

// outside Array class
template <int k, class Array>
struct Array_proxy {

  typedef const Array_proxy<k-1, Array> value_type;
  typedef Array_proxy<k-1, Array> reference_type;

  explicit Array_proxy (const Array& a, size_t i) : a_(a), i_(i) {}

  template <int c>
  Array_proxy (const Array_proxy<c, Array>& a, size_t i) : a_(a.a_)
      , i_(a.i_) {

    size_t s = 1;
    for (int j=0; j<Array::k_ - k - 1; ++j)
```

```
      s *= a_.n_[j];
    i_ += i*s;
  }

  reference_type operator[](size_t i)
  { return reference_type(*this, i); }

  value_type operator[](size_t i) const
  { return value_type(*this, i); }

  const Array& a_;
  size_t i_;
};

template <class Array>
struct Array_proxy<0, Array> {

  typedef typename Array::reference_type reference_type;
  typedef typename Array::value_type value_type;

  explicit Array_proxy (const Array& a, size_t i) : a_(a), i_(i) {}

  template <int k>
  Array_proxy (const Array_proxy<k, Array>& a, size_t i) : a_(a.a_)
      , i_(a.i_) {

    size_t s = 1;
    for (int j=0; j<Array::k_ - 1; ++j)
      s *= a_.n_[j];
    i_ += i*s;
  }

  reference_type operator=(value_type v) {
    a_.data_[i_] = v;
    return a_.data_[i_];
  }
  reference_type operator+=(value_type v) {
    a_.data_[i_] += v;
    return a_.data_[i_];
  }
  // and similarly for the other compound operators -=, *=, /=

  operator reference_type()
  { return a_.data_[i_]; }

  operator value_type() const
  { return a_.data_[i_]; }

  const Array& a_;
  size_t i_;
```

```
};
```

The unary **operator[]** inside the **Array** class returns an object of a type defined in the **Array_proxy** class template. The latter defines this type as a proxy object of a lower rank, and the partial template specialization **Array_proxy**<0, **Array**> finishes the compile-time recursion. There are two implementations of the unary **operator[]**, one for reading from **const Array** objects and another for read/write access. The aforementioned specialized template class also implements operators so that elements of the **Array** can be modified by the usual assignment operator =, and by the compound assignment operators +=, −=, ∗=, /=.

A similar metaprogram can be written to output the array to a stream from the standard library. Here, a class template **Print** implements a print static member function that calls the same function on a class of a lower rank recursively at compilation time, until the recursion is finished with **Print**<2>. In the case of a vector object, the template specialization **Print**<1> handles the printing of the vector elements to the output stream.

The **Array** class template can have a much richer interface, and the metaprogramming techniques described above can be used to define it. If a function is coherent only for an array of specific rank, it can be defined in the corresponding partial template specialization of the **Array_base** class introduced earlier.

## 4   Language expressions

This section layouts the syntax used to work with the **Array** class template presented in the previous section. The methodology presented follows that of expression templates, introduced by Veldhuizen [8]. Now that the **Array** class template is in place, it is desirable to write code as the following:

```
matrix_type<double> A(m,n), B(n,m);
double alpha;
// initialize objects
// ...
matrix_type<double> C = alpha*A*B;
```

Listing 6 presents the overloaded operators that are needed to accomplish the operation above. An alias template is declared for the resulting expression of a scalar-array multiplication. In this way, the use of aliases reduces considerably the amount of code written by the programmer, and makes the code easier to maintain and modify. In the **SAm**<d,T> alias template declaration, **Expr** is a wrapper class template that contains the information about operators and operands without carrying out any computations. As a result, the actual computation can be deferred to a later time when the result of the expression is needed, a technique known as *lazy evaluation*. Lazy evaluation is an important component of creating an efficient syntax for algebraic operations, as the expressions can be modified efficiently so that the least amount of work is done over them to obtain the result. As a

simple example, the expression $\alpha(\beta A)$, with scalars $\alpha, \beta$ and matrix A, would be transformed into $\gamma A, \gamma = \alpha\beta$ before evaluating the expression. Furthermore, the technique can be used to match patterns of complex mathematical expressions so that they can be handled using the higher-level BLAS routines.

**Listing 6: Overloaded operators**

```cpp
// scalar-array multiplication alias template
template <int d, class T>
using SAm = Expr < BinExprOp < ExprLiteral<T>, Array<d,T>, ApMul >
    >;

// operator*(scalar, array)
template <int d, typename S, typename T>
typename enable_if<is_arithmetic<S>::value, SAm <d,T> >::type
operator*(S a, const Array<d,T>& b) {
  typedef typename SAm <d,T>::expression_type ExprT;
  return SAm <d,T>(ExprT(ExprLiteral<T>(a),b));
}

// expression-(scalar-array multiplication) multiplication alias
    template
template <int d, class T, typename A>
using ESAmm = Expr < BinExprOp < Expr <A>, SAm <d,T>, ApMul > >;

// operator*(expr, array)
template<int d, typename T, class A>
ESAmm <d,T,A> operator*(const Expr <A>& a, const Array<d,T>& b) {
    typedef typename ESAmm <d,T,A>::expression_type ExprT;
    return ESAmm <d,T,A>(ExprT(a, T(1)*b));
}
```

The **enable_if** idiom, shown in the definition of **operator\***, was introduced by Järvi et al. [18] to deal with ambiguous function calls in template overload resolution. That is, if the type **S** is not arithmetic, the function template is not considered as a valid candidate for the function call due to the *Substitution Error Is Not An Error (SFINAE)* situation. In other words, if **S** is not an arithmetic type, the compiler removes this function as a candidate instead of signaling a compilation error. Yet, if **S** is indeed an arithmetic type, the function template is a valid candidate for the overload resolution and the return type is defined to **SAm**<d,T>. The return type is constructed by the aid of an **expression_type** type definition declared within the **Expr** class template, explained below in this section.

The second alias template declares the result of multiplying an arbitrary expression by the expression of a scalar-matrix multiplication. It is worth noting that one needs to implement only the operators for the base cases, i.e., those that deal with non-expression objects, as a simple operator taking two expressions automatically generates the required code. Listing 7 shows the definition of **operator\*** taking two expressions.

```
// operator*(expr, expr)
template<class A, class B>
Expr < BinExprOp < Expr <A>, Expr <B>, ApMul > >
operator*(const Expr<A>& a, const Expr<B>& b) {

    typedef BinExprOp < Expr <A>, Expr <B>, ApMul > ExprT;
    return Expr <ExprT>(ExprT(a,b));
}
```

Now we turn our attention to the types that appear in Listings 6 and 7. The **ExprLiteral** class template is used to wrap a constant or literal [8]. The class template provides implicit conversion to the type stored. The **Expr** class template is shown in Listing 8. It is used as a wrapper to a binary expression class template and it exists to minimize the number of operator overloads [8]. The class declares type definitions of the wrapped class, of type **expression_type**, and its interface redirects to the latter. The listing also shows the new optional return value syntax introduced by the standard, where the keyword **auto** replaces the return type and the actual return type is placed at the end of the function declaration. The type of the return type is obtained by using the **decltype** operator, which was also introduced by the `C++11` set of requirements to let the compiler infer the type of an expression.

Listing 8: **Expr** class template

```
// Expression wrapper class
template <class A>
class Expr {

    A a_;

public:

    typedef A expression_type;

    Expr() : a_() {}
    Expr(const A& x) : a_(x) {}

    auto left() const -> decltype(a_.left()) { return a_.left(); }
    auto right() const -> decltype(a_.right()) { return a_.right();
        }

    operator decltype(a_())() { return a_(); }

    auto operator()() const -> decltype(a_()) { return a_(); }

    friend inline std::ostream& operator<<(std::ostream& os, const
        Expr<A>& expr);
};
```

11

The **BinExprOp** class template, displayed in Listing 9, encapsulates the actual binary expression. The class uses the **Expr_traits** traits class to obtain the type of the left and right branches of the expression. This is due to the fact that some objects are better stored by reference so no copy is necessary when handling expressions. The **operator_type** is the class that handles the work to be done on the expression and it will be explained in detail later in this section. The class template can also handle unary expressions, for the resulting expression can add an empty type as the right-branch. For example, the expression that declares the transposition of an array takes the form **Expr< BinExprOp<Array<d,T>, EmptyType, ApTr> >**, where **EmptyType** is an empty **struct** whose sole purpose is to fill in where a template parameter is not needed [9].

Listing 9: **BinExprOp** class template

```
// binary expression class template
template <class A, class B, class Op>
class BinExprOp {

    typename Expr_traits<A>::type a_;
    typename Expr_traits<B>::type b_;

public:

    typedef Op operator_type;

    BinExprOp(const A& a, const B& b) : a_(a), b_(b) {}

    auto left() const -> decltype(a_) { return a_; }
    auto right() const -> decltype(b_) { return b_; }

    auto operator()() const -> decltype(Op::apply(a_, b_))
    { return operator_type::apply(a_, b_); }
};
```

The only remaining component to discuss in detail is the operator applicative classes. The applicative multiplication class **ApMul** is given in Listing 10. The class shows the functions needed to execute the code **matrix_type<double>** C = alpha*A*B. An alias template for the multiplication between a scalar and a matrix is declared using the scalar-array alias template of Listing 6. The static function shown first in the listing is used to carry out the multiplication between two matrices, each one with its respective scalar. In other words, this function can be used to obtain C = $(\alpha A)(\beta B)$. In this implementation, the function calls obtains the result of the multiplication by calling the respective BLAS function, but the developer is allowed to call another function or to provide her own implementation. However, since there is a BLAS library that is highly optimized for almost every computer architecture, it is strongly recommended to follow this methodology. More static functions can be needed in order to provide further functionality as the syntax of the expressions evolve.

**Listing 10: ApMul multiplication applicative class**

```cpp
// scalar − matrix multiplication alias template
template <typename T>
using SMm = SAm <2,T>;

// Multiplication applicative class
class ApMul {
public:

    // ... other applicative functions

    // SVm − SVm multiplication
    template <typename T>
    static matrix_type<T> apply (const SMm <T>& x, const SMm <T>& y)
        {

      // get matrix refernces
      const matrix_type<T>& a = x.right();
      const matrix_type<T>& b = y.right();

      // check size
      assert (a.columns() == b.rows());

      matrix_type<T> r(a.rows(), b.columns());
      cblas_gemm(CblasNoTrans, CblasNoTrans, r.rows(), r.columns(),
                 a.columns(), x.left()*y.left(),
                 a.data_, a.rows(), b.data_, b.rows(), 1.0, r.data_,
                    r.rows());
      return r;
    }

    // ... other applicative functions

    // expr − expr multiplication
    template<class A, class B>
    static typename Return_type<Expr<A>, Expr<B>, ApMul>::
        result_type
    apply (const Expr<A>& a, const Expr<B>& b)
    { return a()*b(); }
};
```

A catch-all template function is also needed, so that the multiplication between unknown expressions calls **operator()** recursively until a base case, as the one shown earlier for the multiplication between matrices, finishes the recursion. Listing 11 shows the metaprogram used to obtain the **Return_type** of an expression, as declared by the return type in this function. The first four definitions of the **Return_type** class template are general and deal with the operations of expressions and binary expression objects. When dealing with a complicated expression that encapsulates many more complex expressions, the metapro-

gram uses these general classes to obtain the return type of the wrapped objects recursively at compile time. The recursion ends with the partial template specializations that deal with operations between base objects, e.g., an array or a literal. For example, the return type of an operation between a scalar and a matrix defines the matrix as the return type. This does not mean that we can add a scalar to a matrix, as the applicative class has to define the supported operations. The following partial template specialization in the listing defines the return type for the multiplication between a transposed vector and a vector, i.e., for a scalar product. The final partial template specialization defines the return type for the dyadic product between two vectors. Many other specializations can be implemented depending on the needs of the developer.

**Listing 11: Return_type metaprogram**

```cpp
template <typename... Params>
struct Return_type;

// return type for an arbitrary binary expression with arbitrary
    operation
template <typename A, typename B, class Op>
struct Return_type < Expr < BinExprOp <A, B, Op > > > {
    typedef typename Return_type<A>::result_type left_result;
    typedef typename Return_type<B>::result_type right_result;
    typedef typename Return_type<left_result, right_result, Op >::
        result_type result_type;
};

// return type for the operation of an arbitrary object with an
    arbitrary expression
template <typename A, typename B, class Op>
struct Return_type < A, Expr < B >, Op > {
    typedef typename Return_type<typename B::left_type,
    typename B::right_type, typename B::operator_type >::result_type
        right_result;
    typedef typename Return_type<A, right_result, Op >::result_type
        result_type;
};

// return type for the operation of an arbitrary expression with an
    arbitrary object
template <typename A, typename B, class Op>
struct Return_type < Expr<A>, B, Op > {
    typedef typename Return_type<typename A::left_type,
    typename A::right_type, typename A::operator_type >::result_type
        left_result;
    typedef typename Return_type<left_result, B, Op >::result_type
        result_type;
};
```

```cpp
// return type for the operation between two arbitrary expressions
template <typename A, typename B, class Op>
struct Return_type < Expr<A>, Expr<B>, Op > {
    typedef typename Return_type<typename A::left_type,
    typename A::right_type, typename A::operator_type >::result_type
        left_result;
    typedef typename Return_type<typename B::left_type,
    typename B::right_type, typename B::operator_type >::result_type
        right_result;
    typedef typename Return_type<left_result, right_result, Op >::
        result_type result_type;
};

// partial template specializations for scalar - matrix operation
template <int d, typename T, class Op>
struct Return_type<ExprLiteral<T>, Array<d, T>, Op> {
    typedef Array<d, T> result_type;
};

// transposed vector - vector multiplication
template <typename T>
struct Return_type<BinExprOp<Array<1,T>, EmptyType, ApTr>, Array<1,
    T>, ApMul> {
    typedef T result_type;
};

// vector - transposed vector multiplication
template <typename T>
struct Return_type<Array<1,T>, BinExprOp<Array<1,T>, EmptyType,
    ApTr>, ApMul> {
    typedef Array<2,T> result_type;
};
// ... further partial template specializations
```

There are a few caveats that need to be explored. The applicative function described returns a matrix object by value. However, the body of the function is simple enough so that the compiler can carry out the Return Value Optimization (RVO). In cases where the RVO is absent, the notion of *move constructors* can be used to avoid the copy of the returned matrix. The case considered fits perfectly for the implementation of move constructors, which are also part of the new C++ standard, as the entire matrix data is allocated in the heap. Nevertheless, it is often required the modification of an existing object with the result of an expression, instead of the creation of a new one. In C++ this is traditionally accomplished on built-in types by using the compound assignment operators, i.e., operators +=, −=, ∗=, /=. Listing 12 shows the code needed to modify a matrix with the result of an expression that involves scalar and matrix multiplications, i.e., to execute the code C += (alpha*A)(beta*B). It is worth noting that there is no running time performance loss as all the expression template facility is bound to compilation time.

```cpp
// operator+=(any, any)
template <class A, class B>
typename
enable_if<!is_arithmetic<A>::value && !is_arithmetic<B>::value, A&
    >::type
operator+=(A& a, const B& b) {
    typedef RefBinExprOp<A, B, ApAdd> ExprT;
    return Expr<ExprT>(ExprT(a,b))();
}

// (scalar - matrix multiplication) - (scalar - matrix
    multiplication) multiplication alias template1
template <typename T>
using SMmSMmm = Expr < BinExprOp < SMm <T>, SMm <T>, ApMul > >;

// Addition applicative class
class ApAdd {
public:

    // ... other applicative functions

    // array - (scalar*array - scalar*array multiplication)
        addition
    template <typename T>
    static matrix_type<T>& apply(matrix_type<T>& c, const SMmSMmm <
        T>& y) {

        // get matrix refernces
        const matrix_type<T>& a = y.left().right();
        const matrix_type<T>& b = y.right().right();

        // check size
        assert(a.columns() == b.rows());
        assert(c.rows() == a.rows());
        assert(c.columns() == b.columns());

        cblas_gemm(CblasNoTrans, CblasNoTrans, a.rows(), b.columns
            (),
                a.columns(), y.left().left()*y.right().left(),
                a.data_, a.rows(), b.data_, b.rows(), 1.0, c.data_,
                    c.rows());
        return c;
    }
    // other applicative functions
};
```

16

# 5  Usage and performance

The following listing shows an example of the mathematical syntax outlined in the previous section. Clearly, the code is easier to read than writing calls to the BLAS routines.

```cpp
Listing 13: Example usage

#include <iostream>
#include "expr.hpp"

int main() {

    int m, n;
    vector_type<double> x(n), y(n);
    matrix_type<double> A(m,n), B(m,n), C(n,n);
    double alpha, beta;

    // initialize objects
    // ...
    C = x*transpose(y);
    C += alpha*transpose(A)*beta*B;

    std::cout<<"C: "<<C<<std::endl;

    return 0;
}
```

Performance tests were conducted carrying out only matrix-matrix multiplication, i.e., executing the code **matrix_type**<**double**> C = A*B, with square matrix sizes ranging from size $n = 250$ to $n = 8000$. The compiler used was GNU GCC version 4.7, and the program was executed on an Apple MacBook Pro with a 2.66 GHz Intel Core i7 processor. The program links to the BLAS library within the Accelerate.framework for Mac OS X. For each matrix size, the reported time is the result of the average of ten different simulations. Figure 1 shows the results of comparing a direct call to the level 3 BLAS function in Listing 1, with the results of using the **Array** class template described in Section 3 with the mathematical syntax outlined in Section 4. As the figure shows, no performance is lost by using the presented framework for dealing with algebraic objects.

# 6  Conclusion

This article introduces many of the new features provided by the new C++ standard in implementing an efficient building block of any high-performance computing platform. By using operator overloading and expression templates to provide an easy to use syntax and to defer the computation to a later time where the result of an expression is needed, an extra high-level layer is given to the C++ programming language when dealing with
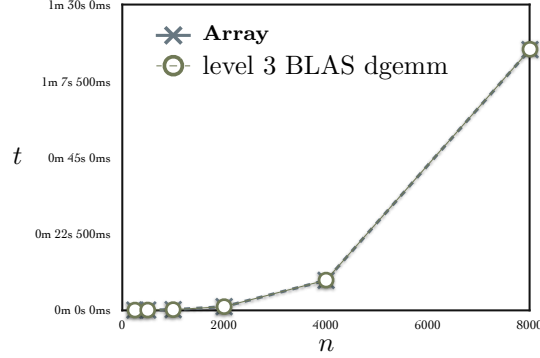
Figure 1: Execution time $t$ for a number of $k$ matrix-matrix multiplication of varying size $n$. The two curves compare the approach described using expression templates and the **Array** class described in the previous sections, with a direct call to the level 3 BLAS function shown in Listing 1.

algebraic objects. Thus, mathematical expressions can be written hiding the details of how operations are implemented without loosing performance.

# References

[1] ISO/IEC IS 9899:1999 Programming languages – C.

[2] ISO/IEC IS 14882:1998 Programming languages – C++.

[3] S. Balay, W. D. Gropp, L. C. McInnes, B. F. Smith, Efficient management of parallelism in object oriented numerical software libraries, in: E. Arge, A. M. Bruaset, H. P. Langtangen (Eds.), Modern Software Tools in Scientific Computing, Birkhäuser Press, 1997, pp. 163–202.

[4] T. Veldhuizen, Arrays in Blitz++, Computing in Object-Oriented Parallel Environments (1998) 501–501.
URL http://dx.doi.org/10.1007/3-540-49372-7_24

[5] J. Siek, A. Lumsdaine, The matrix template library: A generic programming approach to high performance numerical linear algebra, in: D. Caromel, R. Oldehoeft, M. Tholburn (Eds.), Computing in Object-Oriented Parallel Environments, Vol. 1505 of Lecture Notes in Computer Science, Springer Berlin / Heidelberg, 1998, pp. 501–501, 10.1007/3-540-49372-7_6.
URL http://dx.doi.org/10.1007/3-540-49372-7_6

[6] P. Gottschling, D. S. Wise, M. D. Adams, Representation-transparent matrix algorithms with scalable performance, in: Proceedings of the 21st. annual international conference on Supercomputing, ICS '07, ACM, New York, NY, USA, 2007, pp. 116–125.
URL http://doi.acm.org/10.1145/1274971.1274989

[7] A. M. Aragón, cpp-array: A C++ interface to the blas library using arbitrary-rank arrays. (2011).
URL http://code.google.com/p/cpp-array/

[8] T. Veldhuizen, Expression templates, C++ Report 7 (26–31).

[9] A. Alexandrescu, Modern C++ Design: Generic Programming and Design Patterns Applied, Addison-Wesley Professional, 2001.

[10] ISO/IEC IS 14882:2011 Information technology – Programming languages – C++.

[11] D. Gregor, A brief introduction to variadic templates, no. N2087=06-0157, ANSI/ISO C++ Standard Committee Pre-Portland mailing, 2006.

[12] D. Gregor, J. Järvi, J. Maurer, J. Merrill, Proposed Wording for Variadic Templates (Revision 2), no. N2242=07-0102, ANSI/ISO C++ Standard Committee Pre-Portland mailing, 2007.

[13] C. L. Lawson, R. J. Hanson, D. R. Kincaid, F. T. Krogh, Basic linear algebra subprograms for fortran usage, ACM Trans. Math. Softw. 5 (1979) 308–323.
URL http://doi.acm.org/10.1145/355841.355847

[14] J. J. Dongarra, J. Du Croz, S. Hammarling, R. J. Hanson, An extended set of fortran basic linear algebra subprograms, ACM Trans. Math. Softw. 14 (1988) 1–17.
URL http://doi.acm.org/10.1145/42288.42291

[15] J. J. Dongarra, J. Du Croz, S. Hammarling, I. S. Duff, A set of level 3 basic linear algebra subprograms, ACM Trans. Math. Softw. 16 (1990) 1–17.
URL http://doi.acm.org/10.1145/77626.79170

[16] J. Walter, M. Koch. uBLAS – boost basic linear algebra [online].

[17] S. Meyers, More Effective C++: 35 New Ways to Improve Your Programs and Designs, Addison-Wesley Professional, 1996.

[18] J. Järvi, J. Willcock, A. Lumsdaine, Concept-Controlled Polymorphism. In Frank Pfennig and Yannis Smaragdakis, editors, Generative Programming and Component Engineering, Vol. 2830 of LNCS, Springer Verlag, 2003, pp. 228–244.